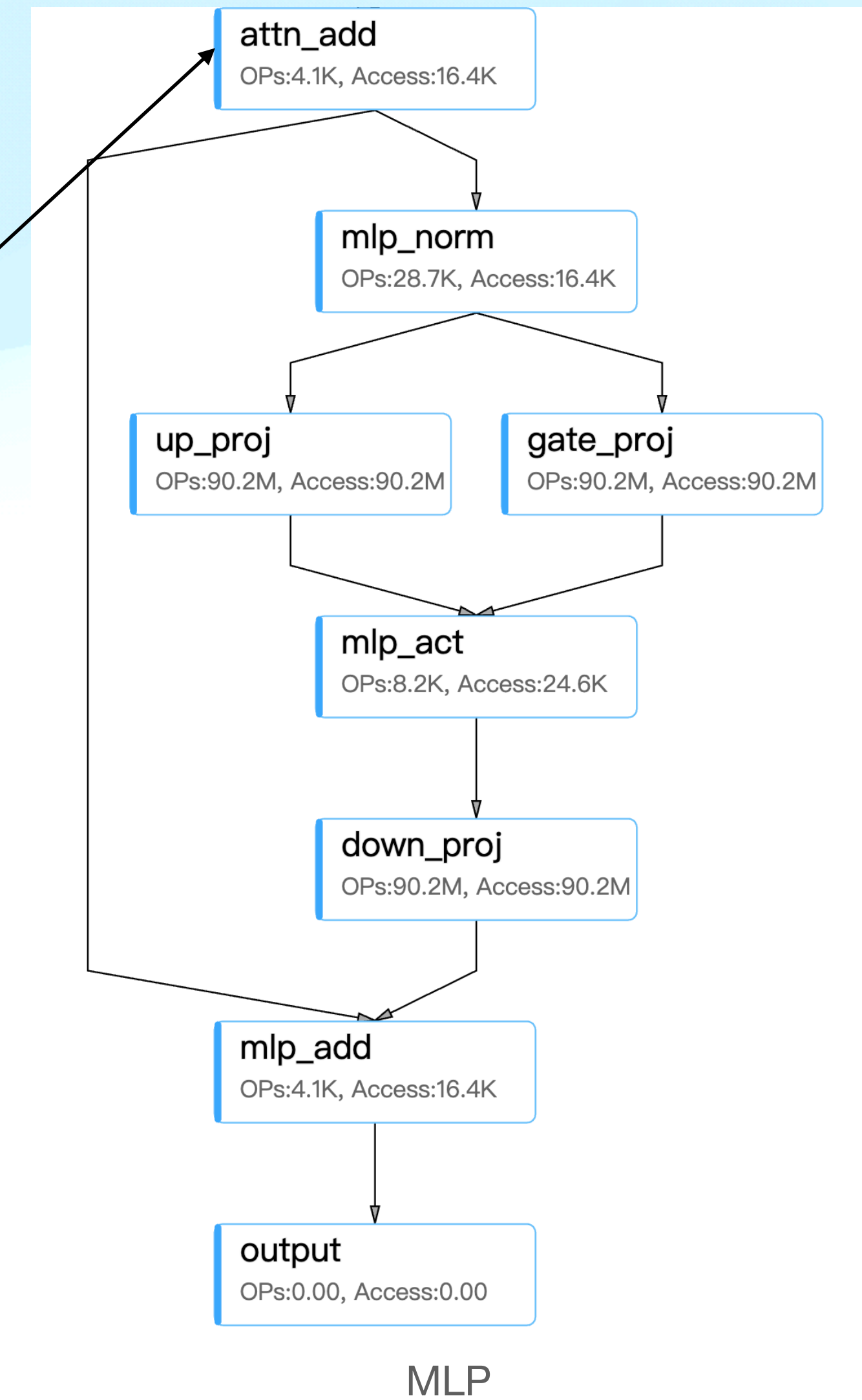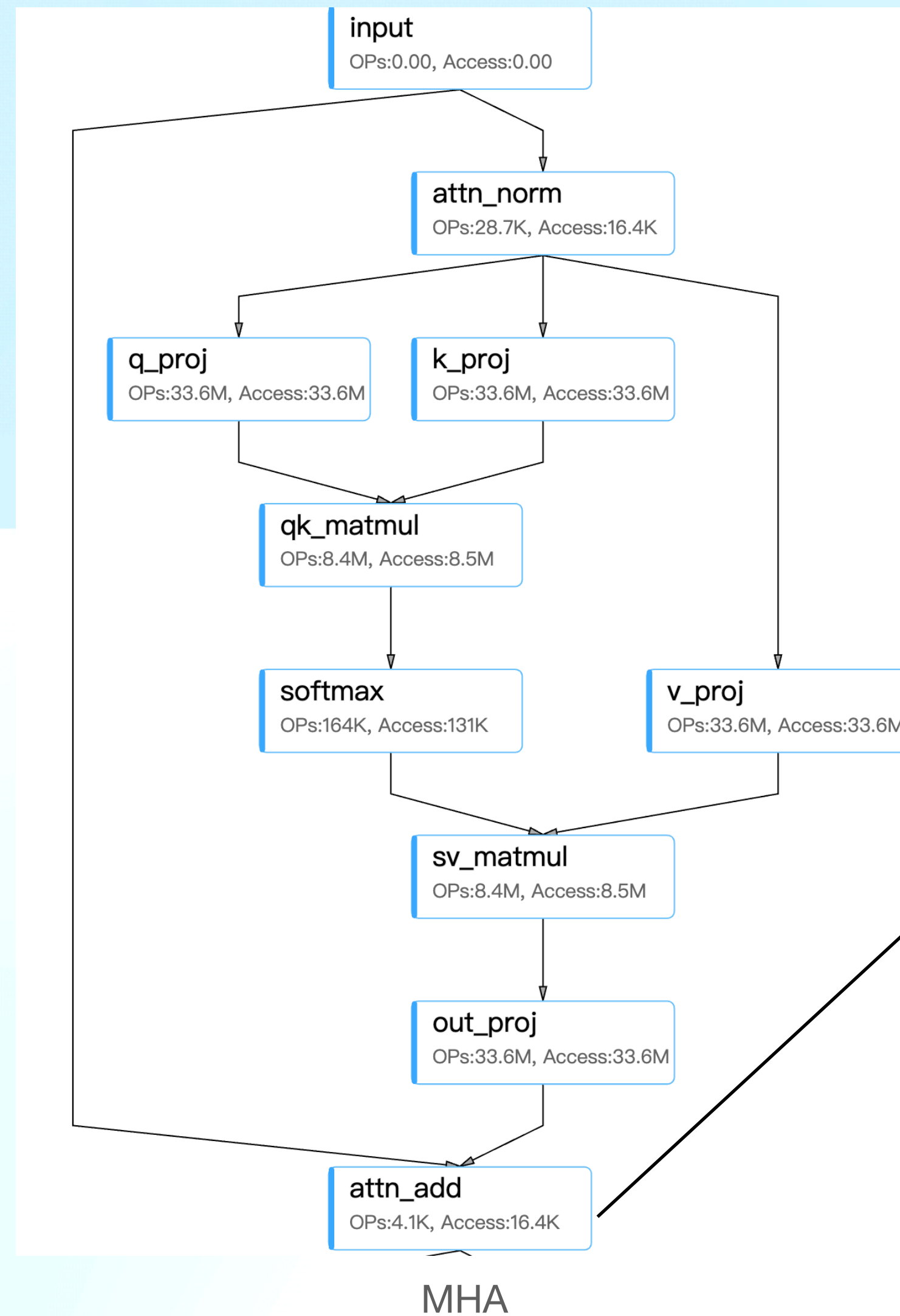# Kernel Optimization

**Profile result**

- Profile

  - End-to-End operation analysis

  - Kernel implemented by Flashinfer

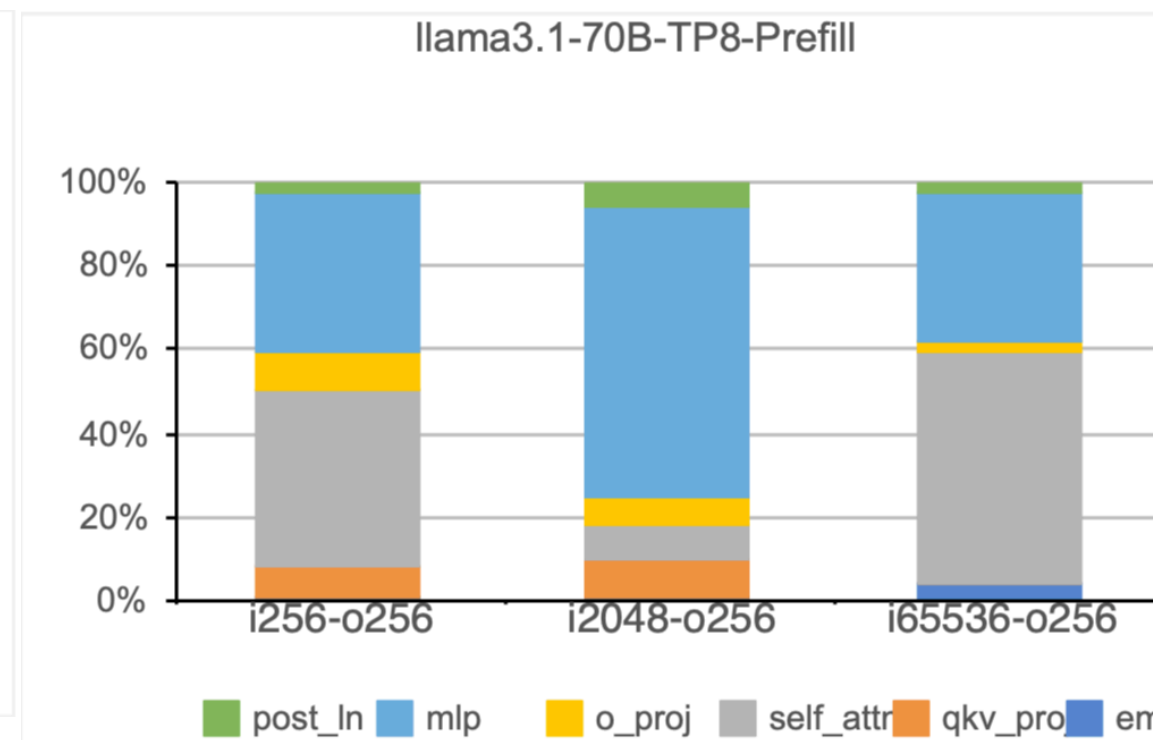- Fuse all kernel to one
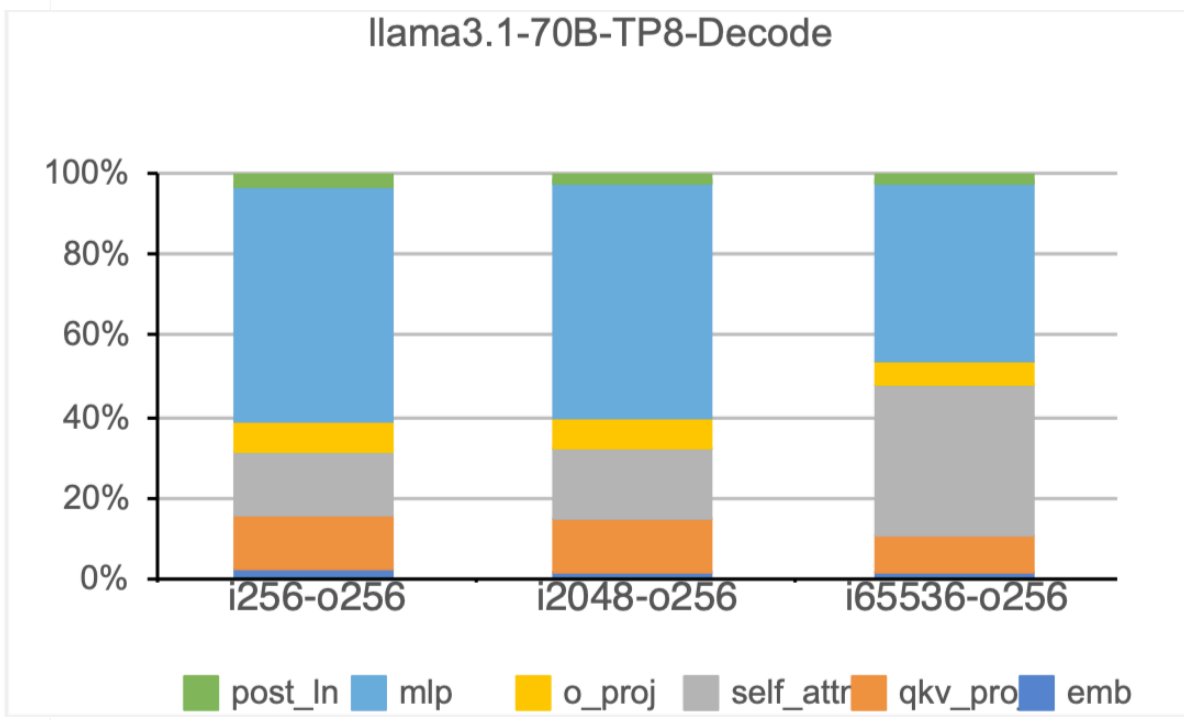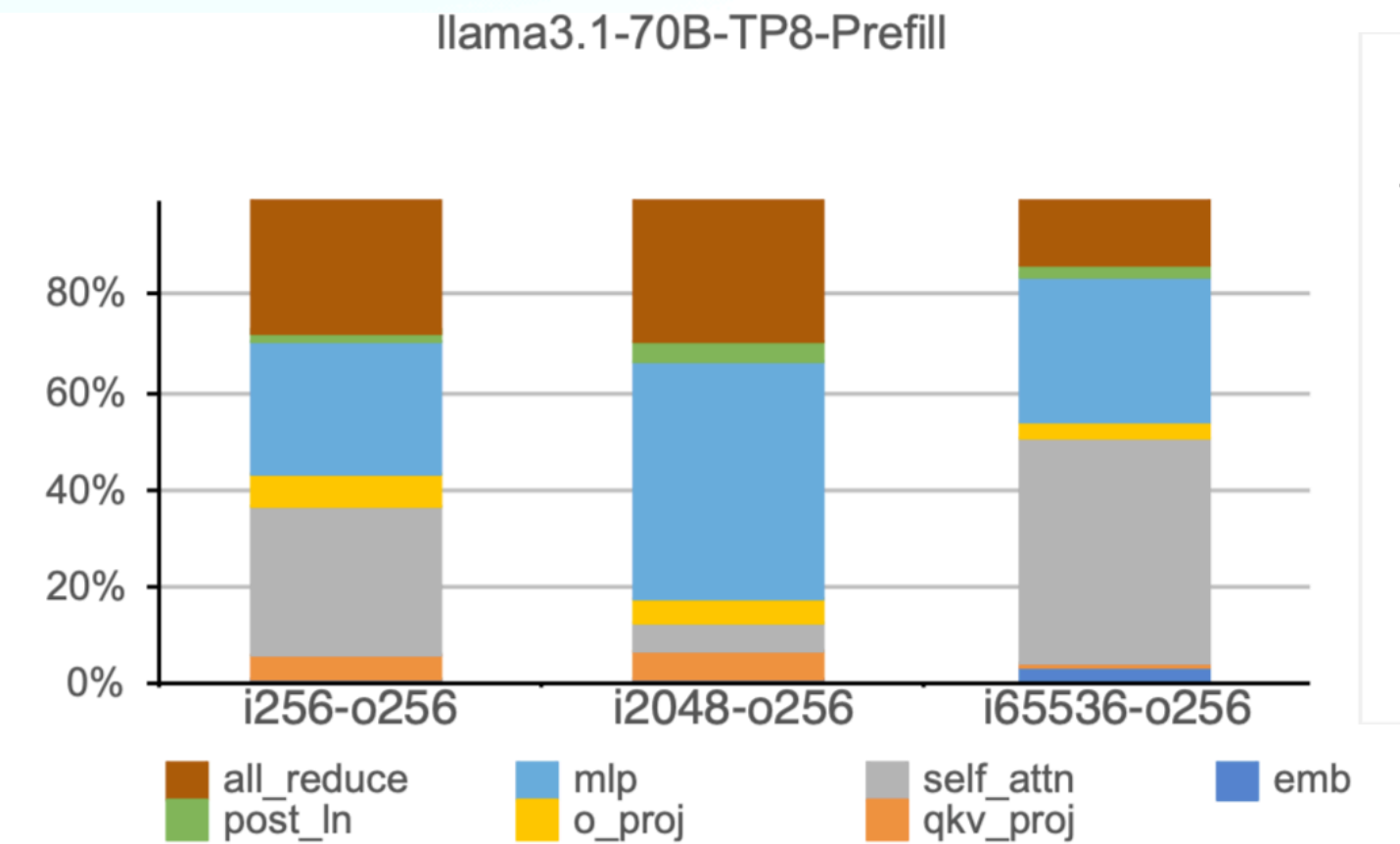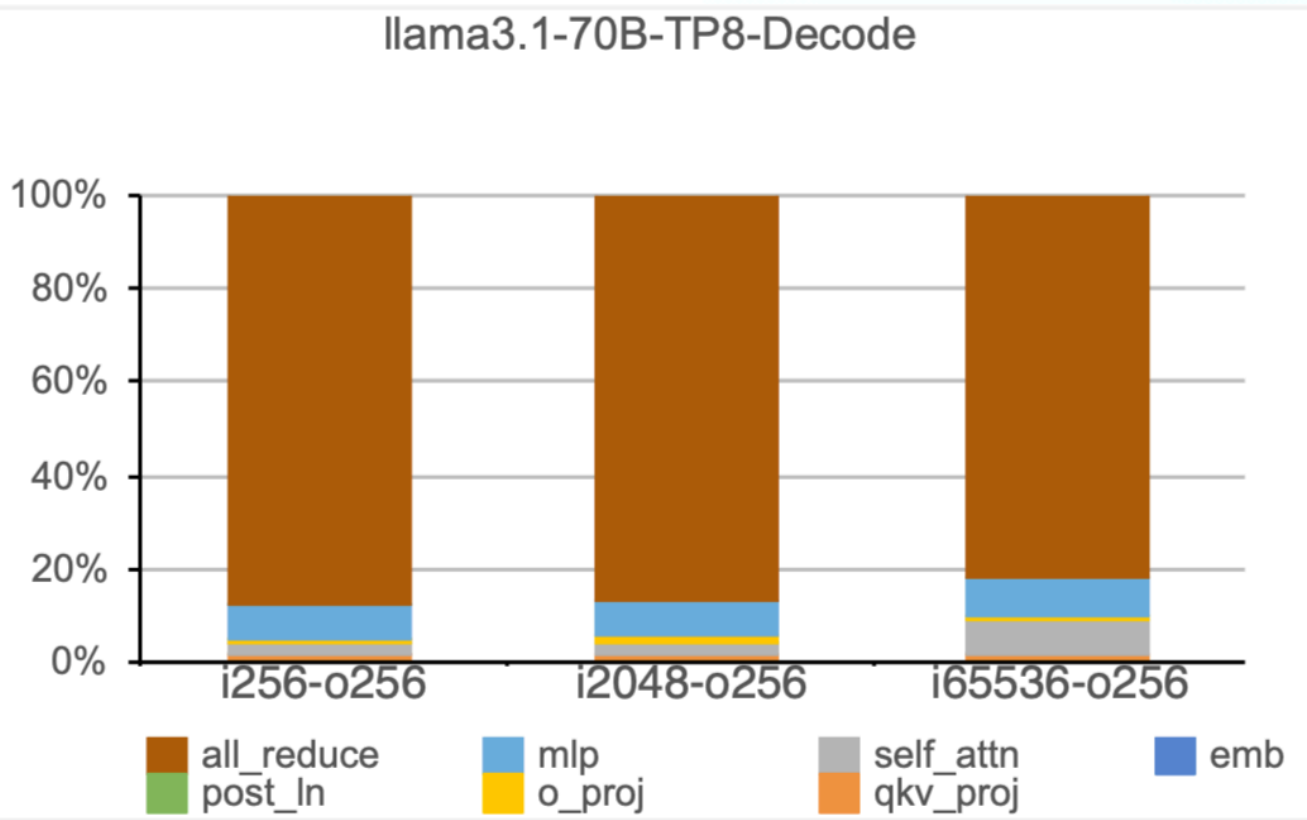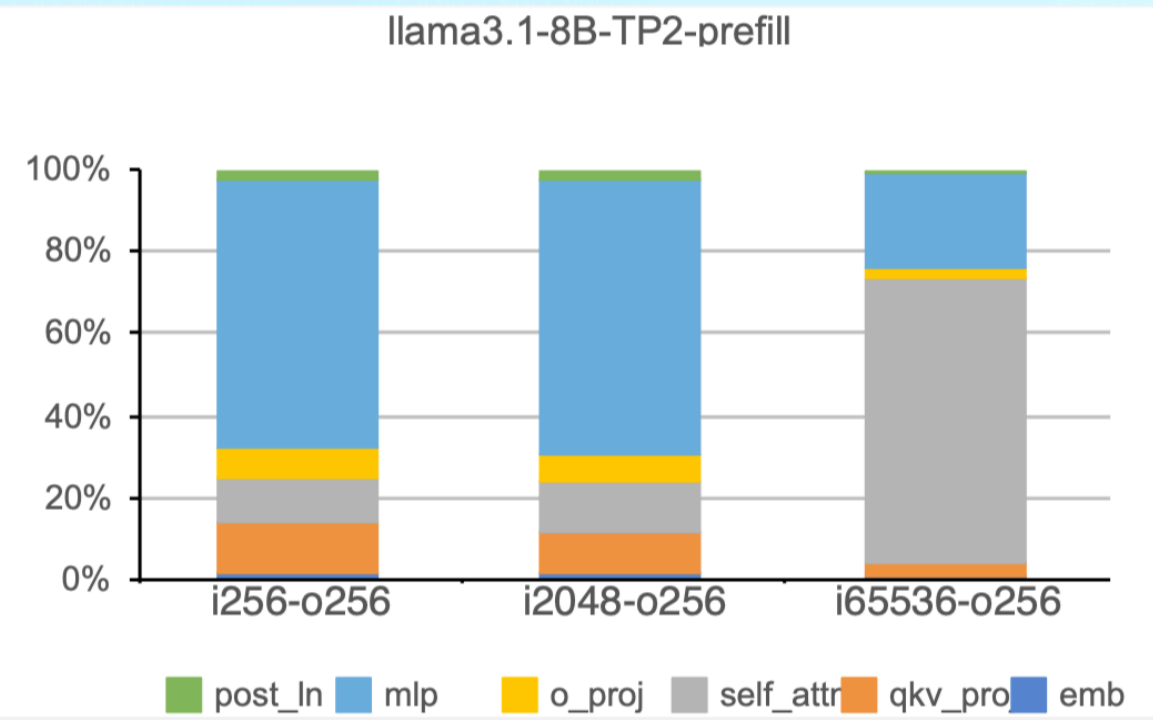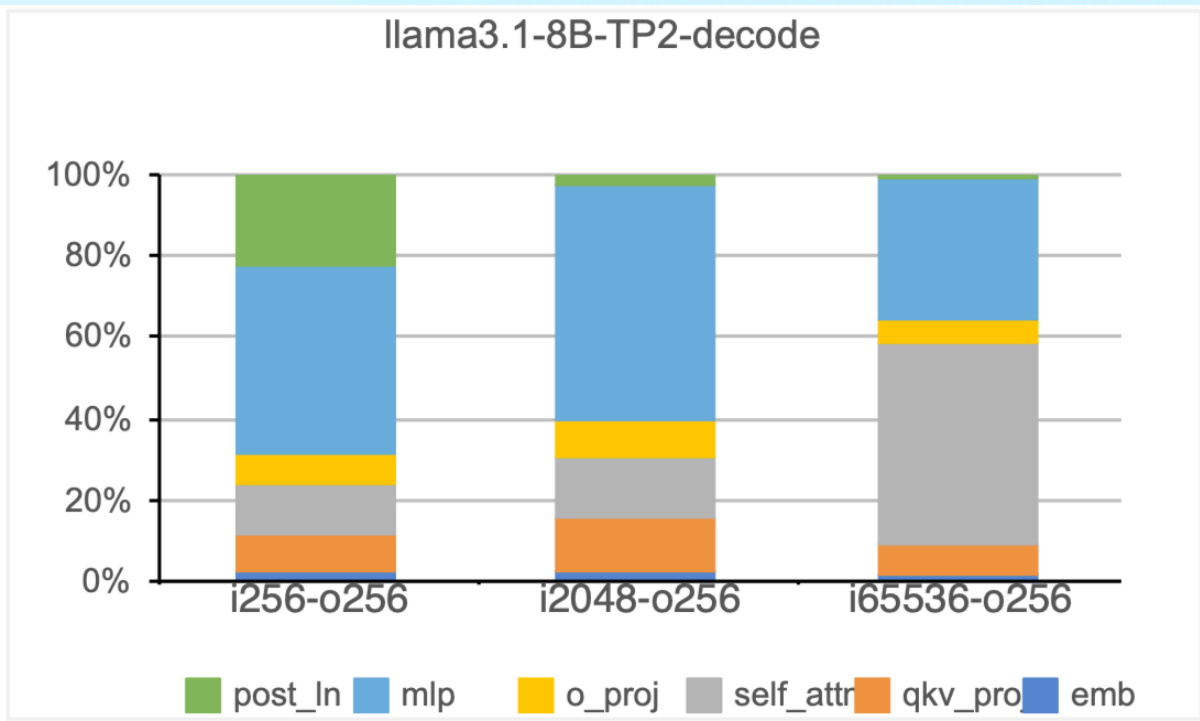
# End-to-End Analysis
## Transformer layer (Llama) operation



**MHA**

- input — OPs:0.00, Access:0.00
- attn_norm — OPs:28.7K, Access:16.4K
- q_proj — OPs:33.6M, Access:33.6M
- k_proj — OPs:33.6M, Access:33.6M
- qk_matmul — OPs:8.4M, Access:8.5M
- softmax — OPs:164K, Access:131K
- v_proj — OPs:33.6M, Access:33.6M
- sv_matmul — OPs:8.4M, Access:8.5M
- out_proj — OPs:33.6M, Access:33.6M
- attn_add — OPs:4.1K, Access:16.4K

**MLP**

- attn_add — OPs:4.1K, Access:16.4K
- mlp_norm — OPs:28.7K, Access:16.4K
- up_proj — OPs:90.2M, Access:90.2M
- gate_proj — OPs:90.2M, Access:90.2M
- mlp_act — OPs:8.2K, Access:24.6K
- down_proj — OPs:90.2M, Access:90.2M
- mlp_add — OPs:4.1K, Access:16.4K
- output — OPs:0.00, Access:0.00

# End-to-End Analysis
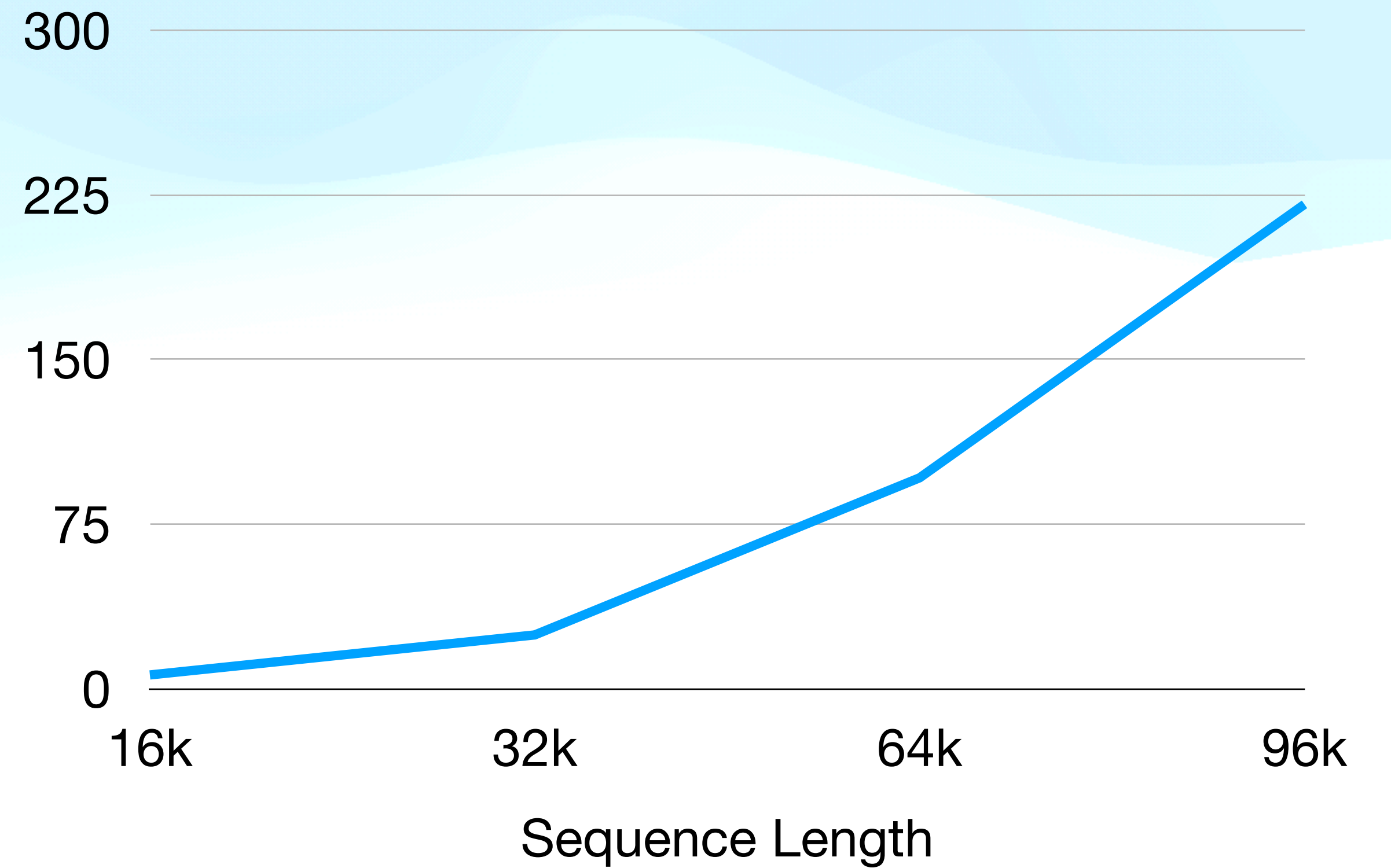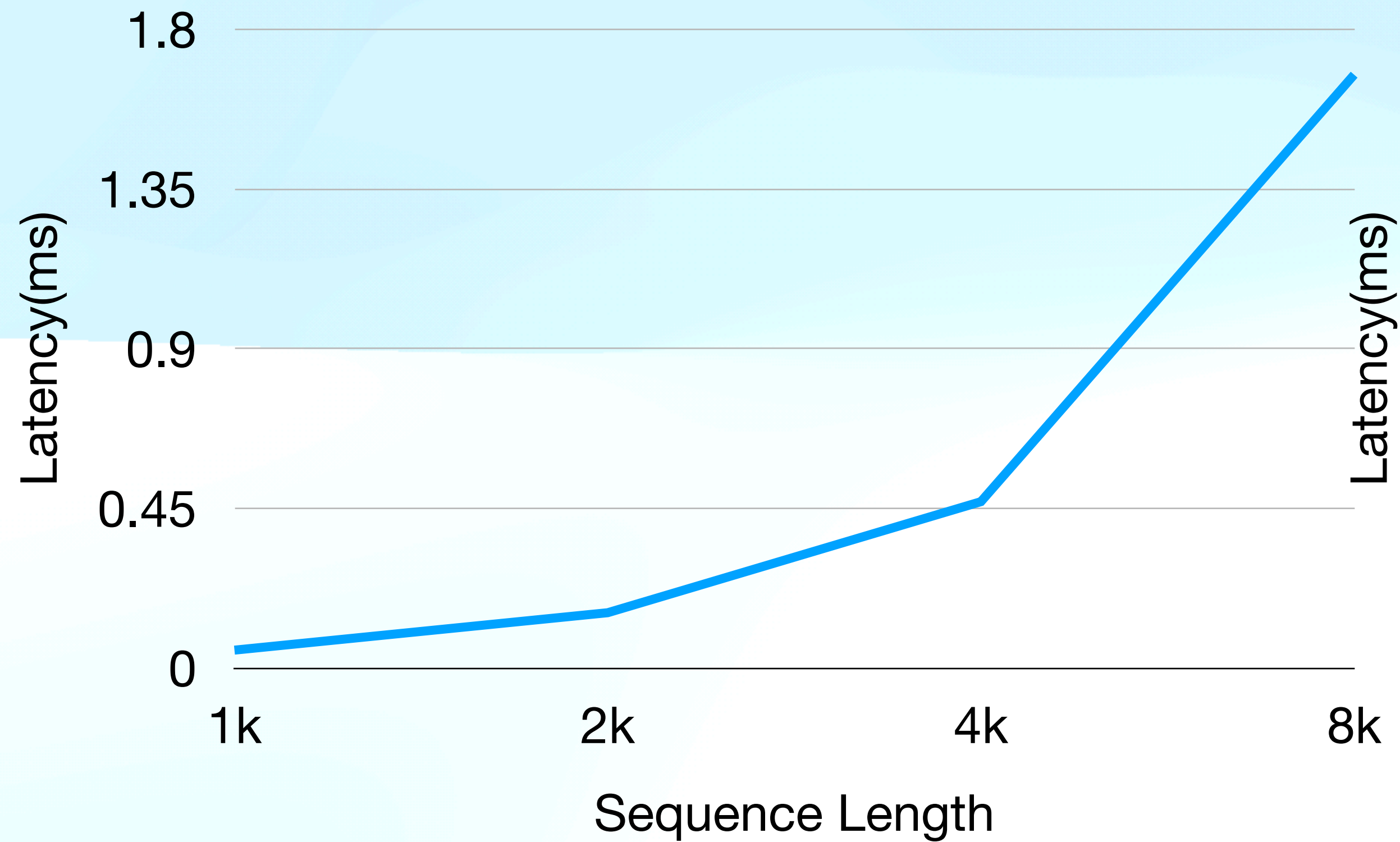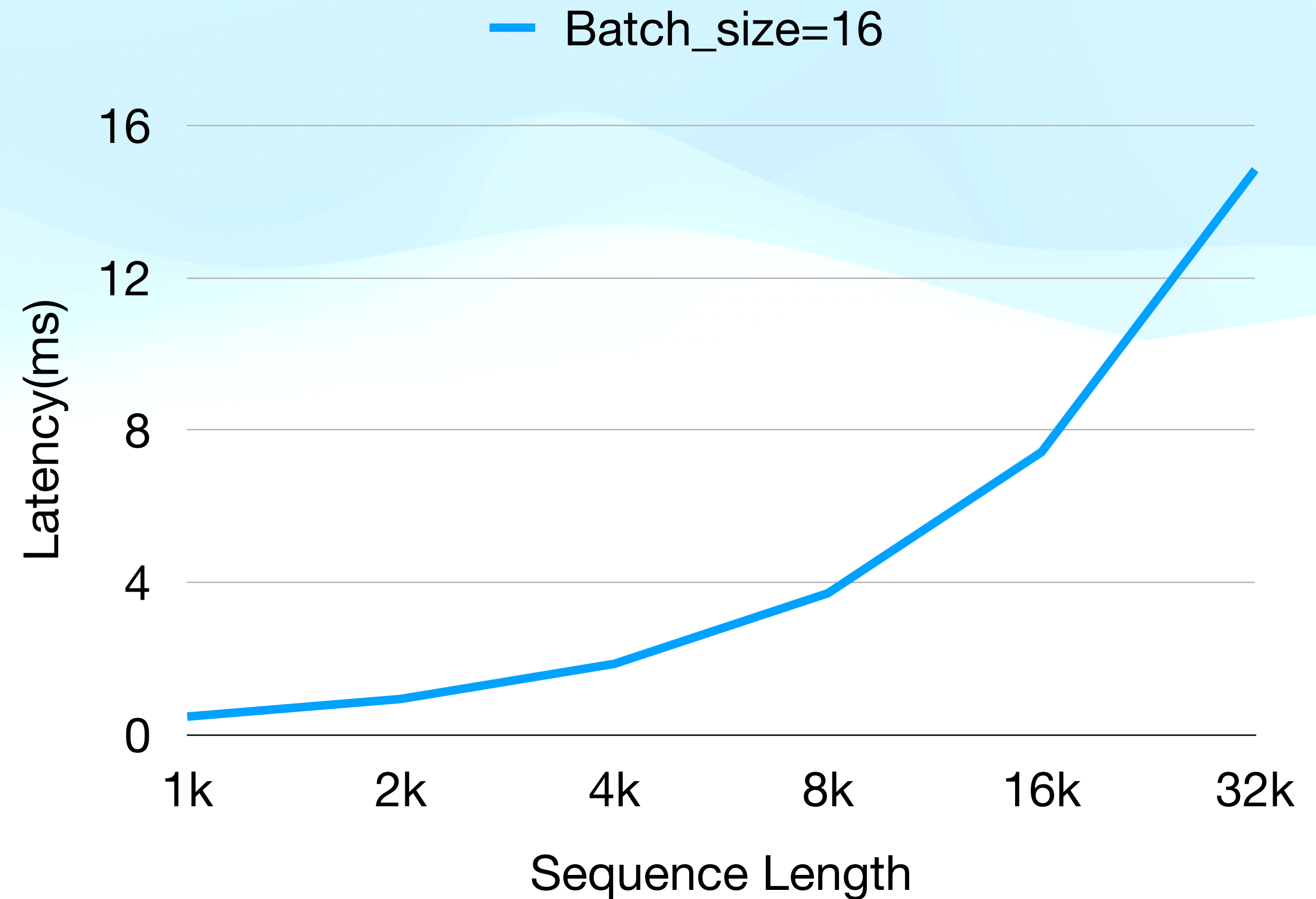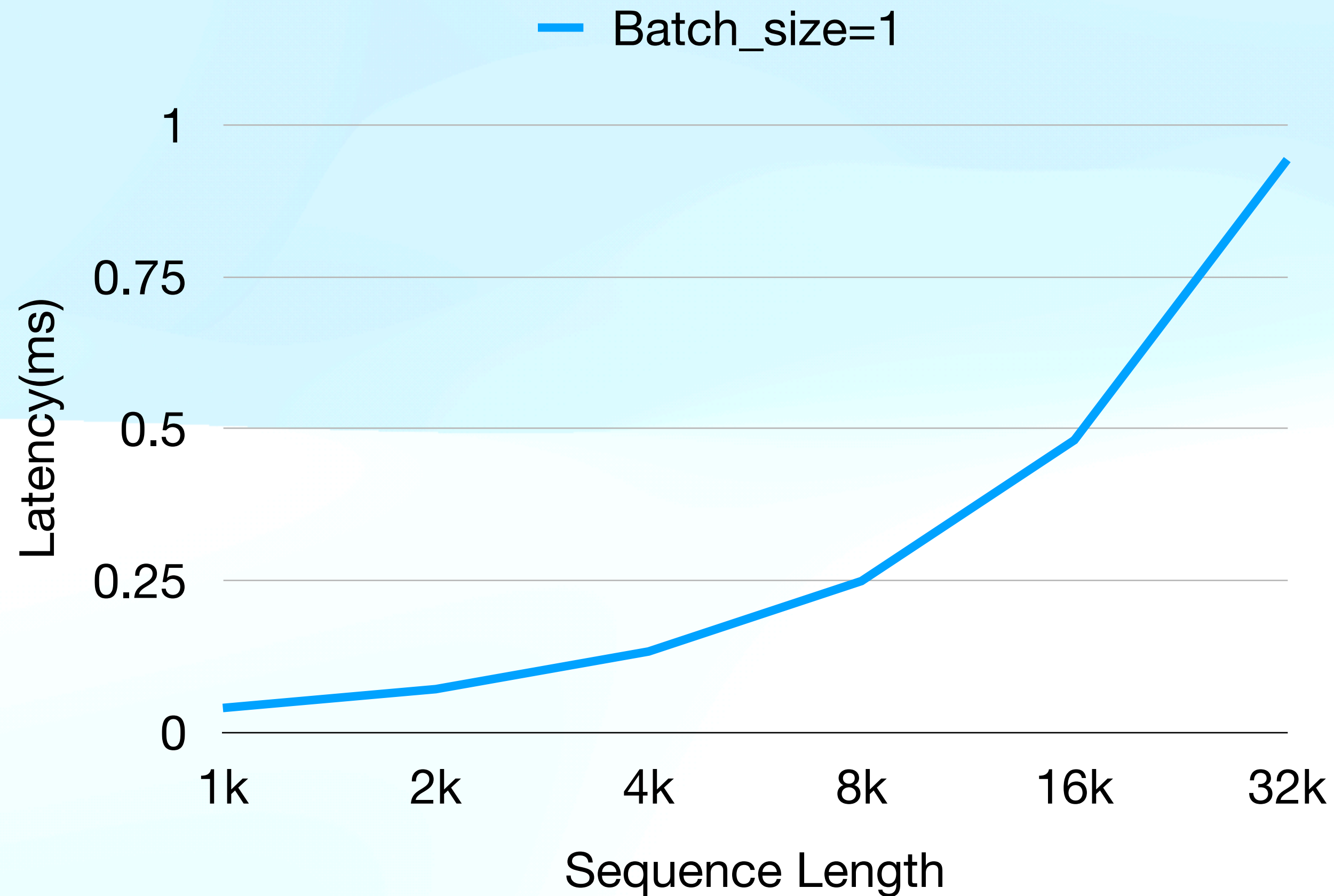## Llama3.1 operation ratio

# Attention Kernel Latency
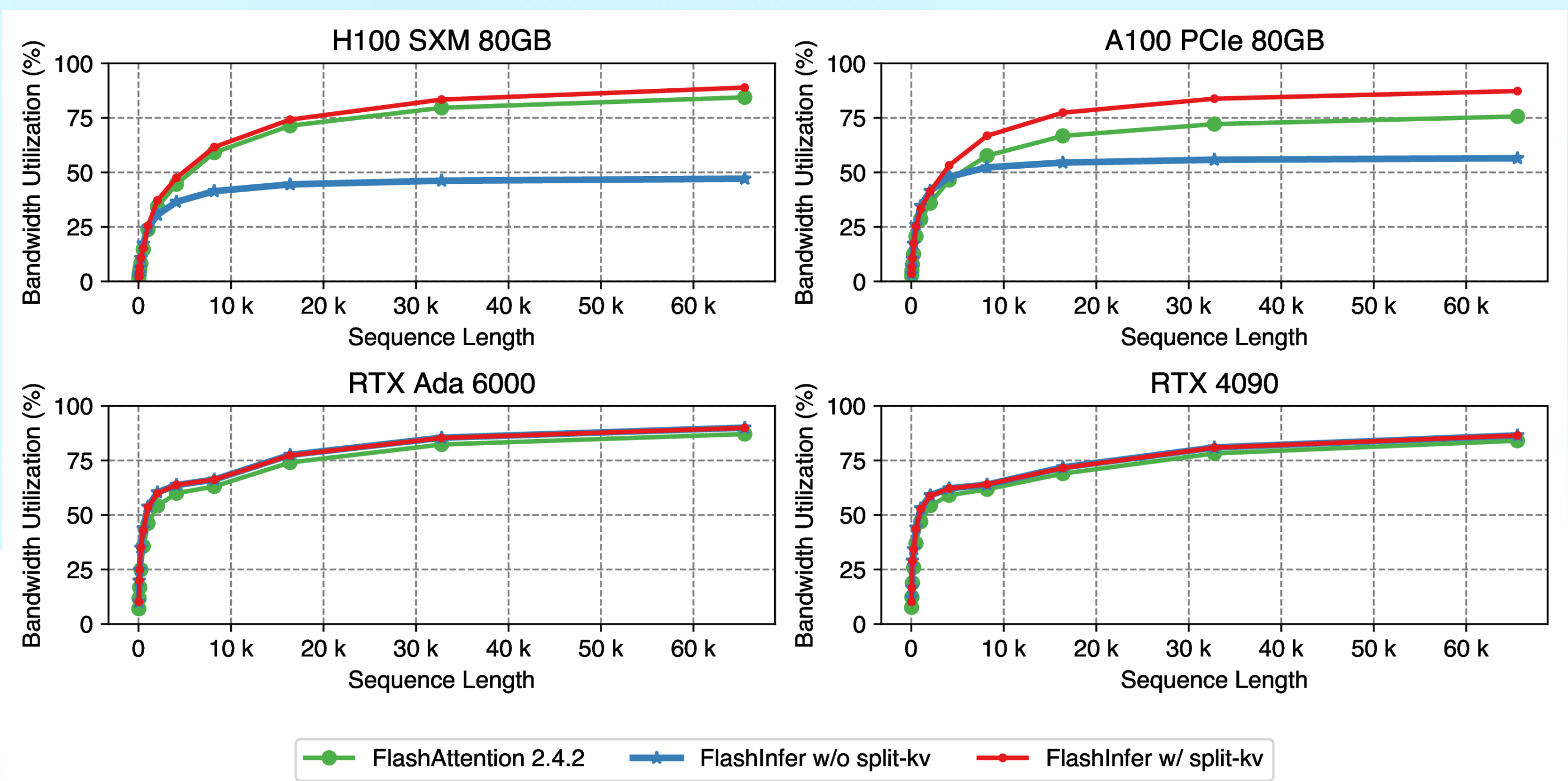
Prefill phase

# Attention Kernel Latency

Decode phase

# Profile Analysis



Single request decode kernel GPU utilization

| | CuBLAS | H100 Peak Tensor Cores Performance |
|---|---|---|
| TFlops | 670 | 989 |

CuBLAS GEMM Performance

We can try to optimize through fusing all kernel to one for better end-to-end performance.
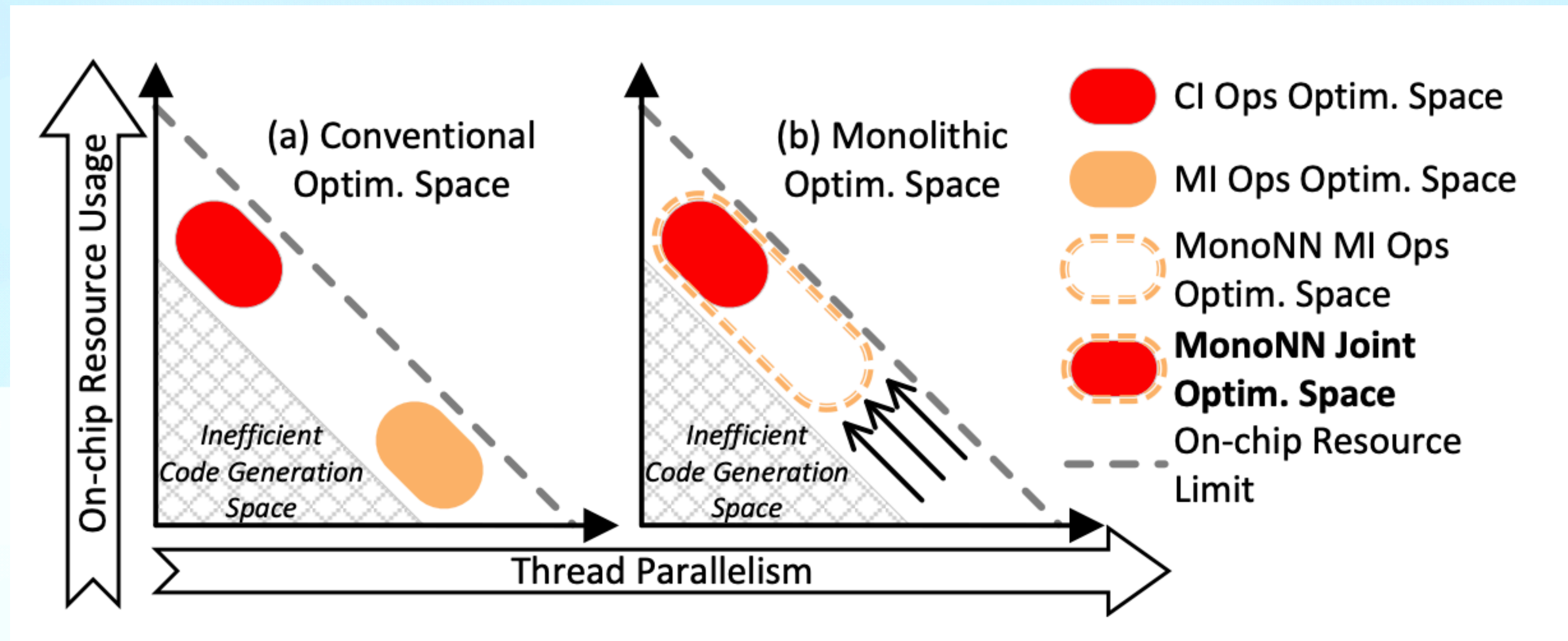
# Fuse all kernel to one
## Challenges

- Resource incompatibility between computeintensive and memory-intensive operators.

    - Compute-intensive operators require a large amount of on-chip resources (registers and shared memory)

    - Memory-intensive ops rely on massive concurrent threads to hide off-chip memory access(TLP).
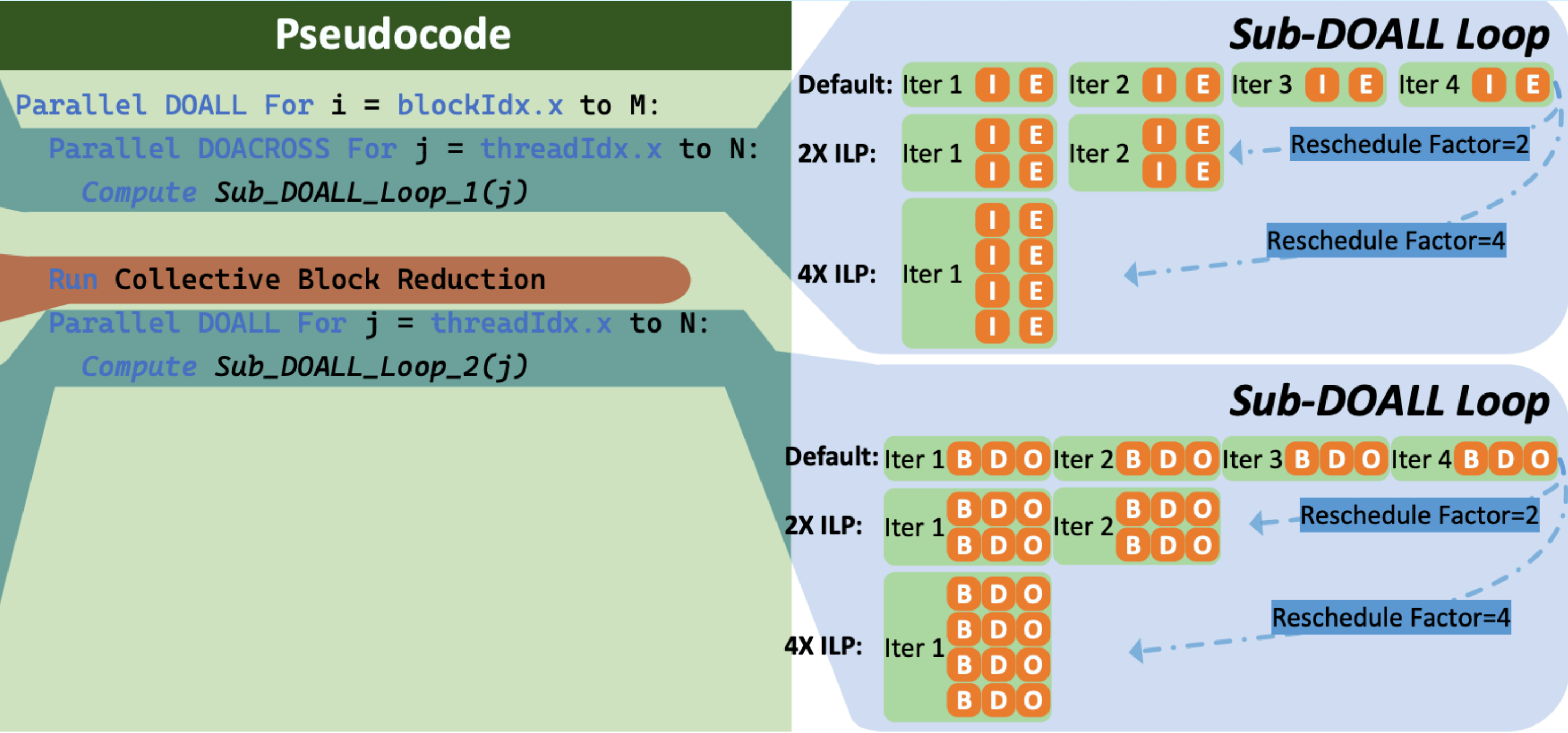
# Fuse all kernel to one
## A High-level illustrate



Align the optimization space of MI Ops as closely as possible with that of CI Ops

# Fuse all kernel to one
## Detailed Implementation



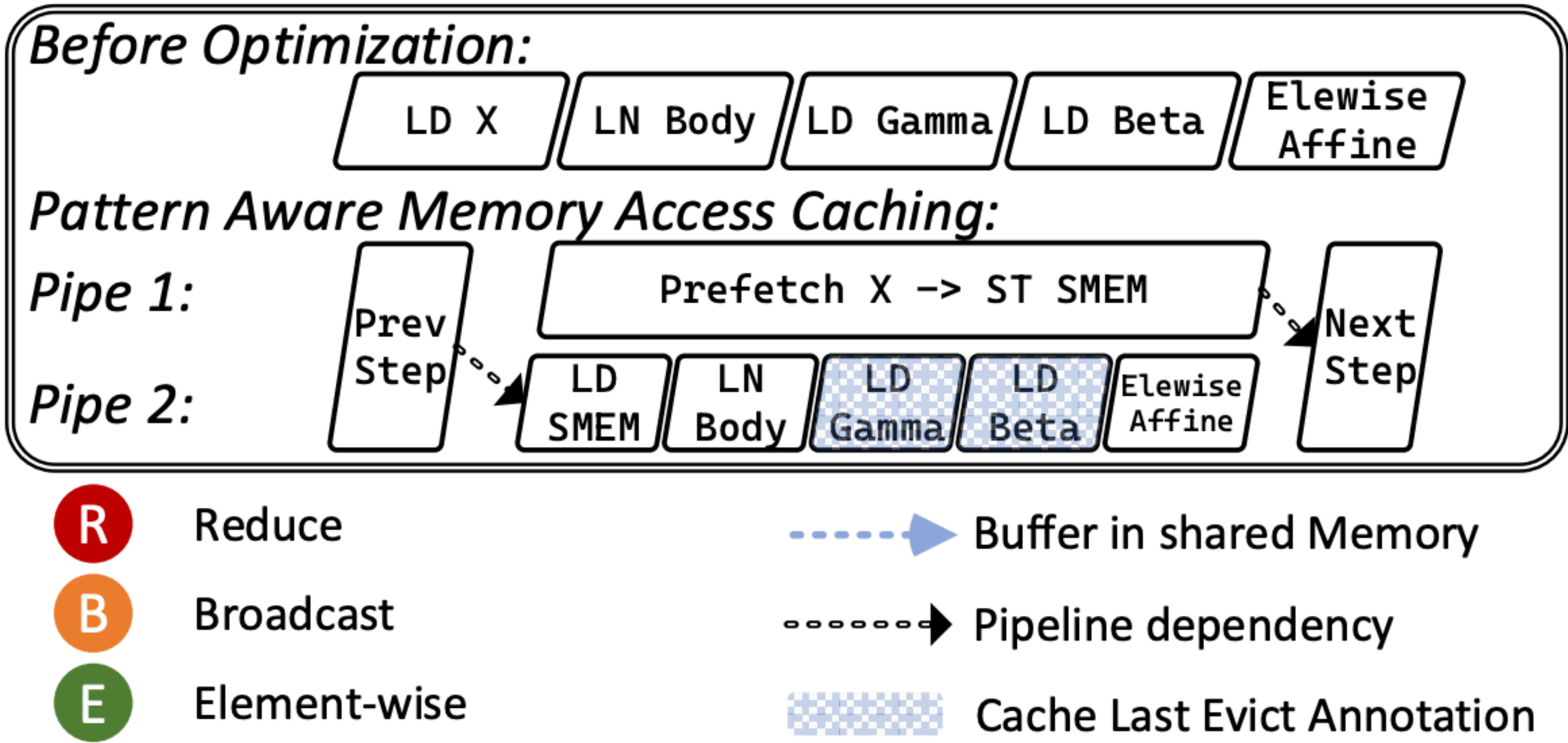Context-aware instruction rescheduling for softmax computation.

Context-aware instruction rescheduling:

For memory-intensive operator(softmax),

we can use ILP(Loop unroll) which need

more on-chip resources adaptive with

compute-intensive operator(GEMM).

# Fuse all kernel to one
## Detailed Implementation



Before Optimization:
LD X | LN Body | LD Gamma | LD Beta | Elewise Affine

Pattern Aware Memory Access Caching:

Pipe 1: Prefetch X -> ST SMEM

Pipe 2: Prev Step | LD SMEM | LN Body | LD Gamma | LD Beta | Elewise Affine | Next Step

R Reduce
B Broadcast
E Element-wise

- - - -> Buffer in shared Memory
· · · · ·> Pipeline dependency
Cache Last Evict Annotation
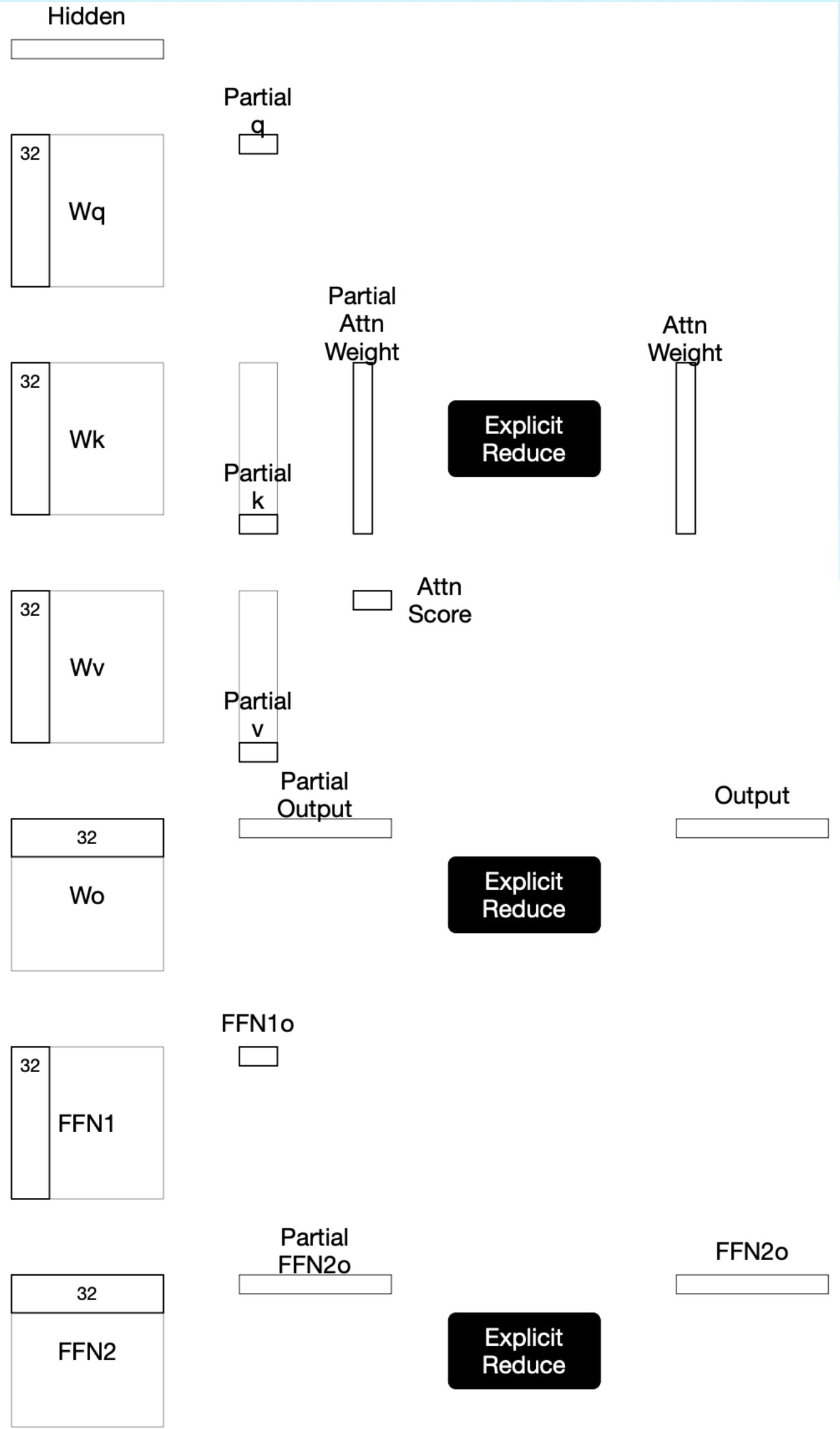
On-chip resource exploitation

Memory access optimizations:

Leverages the abundant shared memory resource

allocated by the compute-intensive computation

in one kernel to pipeline the global memory access

# Fuse all kernel to one

## Fuse decode phase to one

# Thanks